

$X_1, \dots, X_n \sim \mathcal{D}_1$ w/cdf F iid.

$Y_1, \dots, Y_m \sim \mathcal{D}_2$ w/cdf G iid.

Goal: Test $H_0: F=G$ $H_1: F \neq G$

① Estimators

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$$

$$G_m(t) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\}$$

② Pivot / test statistic

$$T_{n,m}$$

$$\psi = \mathbb{1}\{T_{n,m} > s\}$$

③ Adjust s

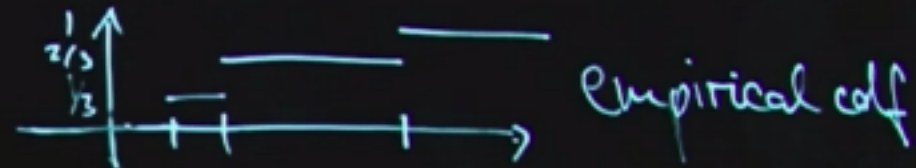
Two-sample Kolmogorov-Smirnov test

Reminder: Kolmogorov-Smirnov goodness of fit test

$X_1, \dots, X_n \sim \mathcal{D}$ w/cdf F iid.

Test $H_0: F = F_0$ $H_1: F \neq F_0$

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$$



$$T_n = \sup_{t \in \mathbb{R}} |F_n(t) - F_0(t)| \xrightarrow[n \rightarrow \infty]{\mathcal{D}} Z \quad \text{supremum of a Brownian bridge}$$

$X_1, \dots, X_n \sim D_1$ w/cdf F
iid.

$Y_1, \dots, Y_m \sim D_2$ w/cdf G
iid.

Assume F continuous &
strictly increasing.

Goal: Test $H_0: F=G$ $H_1: F \neq G$

① Estimators

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$$

$$G_m(t) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\}$$

② Pivot / test statistic

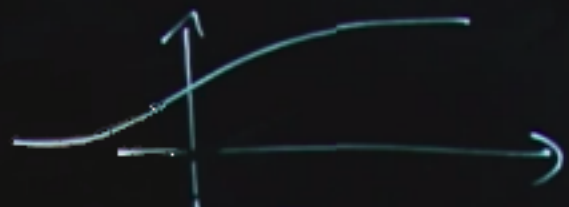
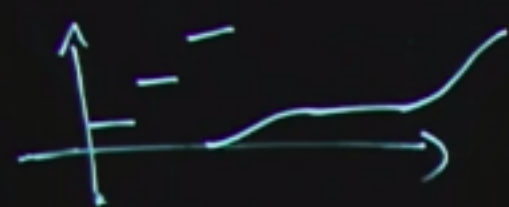
$$T_{n,m} = \sup_{t \in \mathbb{R}} |F_n(t) - G_m(t)|$$

$$\psi = \mathbb{1}\{T_{n,m} > s\}$$

③ Adjust s

Two-sample Kolmogorov-Smirnov test

② Glivenko-Cantelli Thm: $\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow{n \rightarrow \infty} 0$ a.s.



$$\{X_i \leq t\} = \{F(X_i) \leq F(t)\}$$

Asspt $\Rightarrow \exists F^{-1}: (0,1) \rightarrow \mathbb{R}$

$$F^{-1}(F(t)) = t \quad \forall t \in \mathbb{R}$$

$$F(F^{-1}(a)) = a \quad \forall a \in (0,1)$$

$X_1, \dots, X_n \sim D_1$ w/cdf F
 $Y_1, \dots, Y_m \sim D_2$ w/cdf G
iid.

Assume F continuous & strictly increasing.

Goal: Test $H_0: F=G$ $H_1: F \neq G$

① Estimators

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$$

$$G_m(t) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\}$$

② Pivot / test statistic

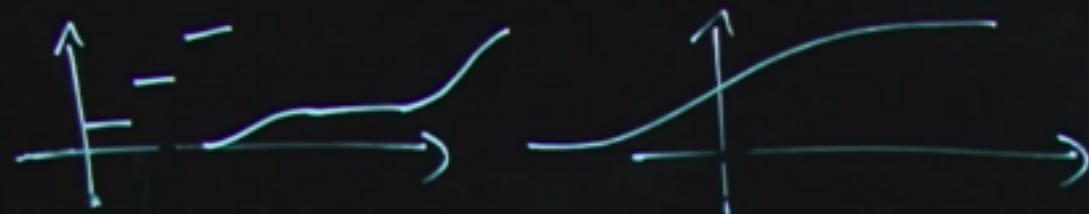
$$T_{n,m} = \sup_{t \in \mathbb{R}} |F_n(t) - G_m(t)|$$

$$\psi = \mathbb{1}\{T_{n,m} > s\}$$

③ Adjust s

Two-sample Kolmogorov-Smirnov test

② "Glivenko-Cantelli-Thm" $\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow{n \rightarrow \infty} 0$ a.s.



Asspt $\Rightarrow \exists F^{-1}: (0,1) \rightarrow \mathbb{R}$

$$F^{-1}(F(t)) = t \quad \forall t \in \mathbb{R}$$

$$F(F^{-1}(a)) = a \quad \forall a \in (0,1)$$

$$\{X_i \leq t\} = \{F(X_i) \leq F(t)\}$$

$$\mathbb{P}(\{F(X_i) \leq a\}) = \mathbb{P}(X_i \leq F^{-1}(a))$$

$$= F(F^{-1}(a)) = a \quad \forall a \in (0,1) \Rightarrow F(X_i) \sim U([0,1])$$

$$T_{n,m} = \sup_t \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} - \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\} \right|$$

$$= \sup_t \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{\underbrace{F(X_i)}_{=U_i \sim U([0,1])} \leq F(t)\} - \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{\underbrace{F(Y_j)}_{=V_j \sim U([0,1])} \leq F(t)\} \right|$$

$$= \sup_{a \in (0,1)} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{U_i \leq a\} - \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{V_j \leq a\} \right|$$

$X_1, \dots, X_n \sim D_1$ w/cdf F
 $Y_1, \dots, Y_m \sim D_2$ w/cdf G
iid.

Assume F continuous &
 strictly increasing.

Goal: Test $H_0: F=G$ $H_1: F \neq G$

① Estimators
 $F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$
 $G_m(t) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\}$

② Pivot/test statistic

$$T_{n,m} = \sup_{t \in \mathbb{R}} |F_n(t) - G_m(t)|$$

$$\psi = \mathbb{1}\{T_{n,m} > s\}$$

③ Adjust s

Two-sample Kolmogorov-Smirnov test

$$T_{n,m} = \sup_{a \in (0,1)} \left| \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{u_i \leq a\} - \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{v_j \leq a\} \right|, u_i, v_j \sim \mathcal{U}(0,1)$$

③ Non-asymptotic: Sample $l=1, \dots, M$ u_i^l, v_j^l ,

$$\text{calculate } T_{n,m}^l = T_{n,m}(u_i^l, v_j^l)$$

$$H_M(q) = \frac{1}{M} \sum_{l=1}^M \mathbb{1}\{T_{n,m}^l \leq q\}$$

Pick $s = \hat{q}_\alpha$, $1-\alpha$ quantile of H_M .

$\approx q_\alpha$, $1-\alpha$ - quantile of cdf of $T_{n,m}$

$$\Rightarrow P(T_{n,m} > s) \approx \alpha$$

$X_1, \dots, X_n \sim \mathcal{D}_1$ w/cdf F
iid.

$Y_1, \dots, Y_m \sim \mathcal{D}_2$ w/cdf G
iid.

Assume F continuous &
strictly increasing.

Goal: Test $H_0: F=G$ $H_1: F \neq G$

① Estimators

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\}$$

$$G_m(t) = \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\}$$

② Pivot / test statistic

$$T_{n,m} = \sup_{t \in \mathbb{R}} |F_n(t) - G_m(t)|$$

$$\psi = \mathbb{1}\{T_{n,m} > s\}$$

③ Adjust s

Two-sample Kolmogorov-Smirnov test

③ Asymptotic KS test: $T_n = \sup_{t \in \mathbb{R}} |F_n(t) - F_0(t)|$

$$\begin{aligned} \text{Fix } t, \text{ Var}(F_n(t)) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{1}\{X_i \leq t\}}_{\sim \text{Be}(F(t))}\right) \\ &= \frac{1}{n^2} \cdot n \cdot F(t)(1-F(t)) = \frac{1}{n} F(t)(1-F(t)) \end{aligned}$$

$$\begin{aligned} \text{Now, for } T_{n,m}: \text{Var}(F_n(t) - G_m(t)) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{X_i \leq t\} - \frac{1}{m} \sum_{j=1}^m \mathbb{1}\{Y_j \leq t\}\right) \\ &= \frac{1}{n^2} \cdot n F(t)(1-F(t)) + \frac{1}{m^2} \cdot m F(t)(1-F(t)) \\ &= \left(\frac{1}{n} + \frac{1}{m}\right) F(t) \cdot (1-F(t)) \end{aligned}$$

$$\frac{n+m}{n \cdot m} \Rightarrow \tilde{T}_{n,m} = \sqrt{\frac{n \cdot m}{n+m}} \cdot \sup_{t \in \mathbb{R}} |F_n(t) - G_m(t)|$$

$\xrightarrow[n, m \rightarrow \infty]{D} Z$, sup of Brownian bridge